

Anomaly Detection and Localization in Crowded Scenes

Weixin Li, *Student Member, IEEE*, Vijay Mahadevan, *Member, IEEE*, and Nuno Vasconcelos, *Senior Member, IEEE*

Abstract—The detection and localization of anomalous behaviors in crowded scenes is considered, and a joint detector of temporal and spatial anomalies is proposed. The proposed detector is based on a video representation that accounts for both appearance and dynamics, using a set of mixture of dynamic textures models. These models are used to implement 1) a center-surround discriminant saliency detector that produces spatial saliency scores, and 2) a model of normal behavior that is learned from training data and produces temporal saliency scores. Spatial and temporal anomaly maps are then defined at multiple spatial scales, by considering the scores of these operators at progressively larger regions of support. The multiscale scores act as potentials of a conditional random field that guarantees global consistency of the anomaly judgments. A data set of densely crowded pedestrian walkways is introduced and used to evaluate the proposed anomaly detector. Experiments on this and other data sets show that the latter achieves state-of-the-art anomaly detection results.

Index Terms—Video analysis, surveillance, anomaly detection, crowded scene, dynamic texture, center-surround saliency

1 INTRODUCTION

SURVEILLANCE video is extremely tedious to monitor when events that require follow-up have very low probability. For crowded scenes, this difficulty is compounded by the complexity of normal crowd behaviors. This has motivated a surge of interest in anomaly detection in computer vision [1], [2], [3], [4], [5], [6], [7], [8], [9]. However, this effort is hampered by general difficulties of the anomaly detection problem [10]. One fundamental limitation is the lack of a universal definition of anomaly. For crowds, it is also infeasible to enumerate the set of anomalies that are possible in a given surveillance scenario. This is compounded by the sparseness, rarity, and discontinuity of anomalous events, which limit the number of examples available to train an anomaly detection system.

One common solution to these problems is to define anomalies as events of low probability with respect to a probabilistic model of normal behavior. This enables a statistical treatment of anomaly detection, which conforms with the intuition of anomalies as events that deviate from the expected [10]. However, it introduces a number of challenges. First, it makes anomalies dependent on the *scale* at which normalcy is defined. A normal behavior at a fine visual scale may be perceived as highly anomalous when a larger scale is considered, or vice versa. Hence, normalcy

models must be defined at multiple scales. Second, different tasks may require *different models of normalcy*. For instance, a detector of freeway speed limit violations will rely on normalcy models based on speed features. On the other hand, appearance is more important for the detection of carpool lane violators, i.e., single-passenger vehicles in carpool lanes. Third, crowded scenes require normalcy models robust to *complex scene dynamics*, involving many independently moving objects that occlude each other in complex ways, and can have low resolution.

In result, anomaly detection can be extremely challenging. While this has motivated a great diversity of solutions, it is usually quite difficult to objectively compare different methods. Typically, these combine different representations of motion and appearance with different graphical models of normalcy, which are usually tailored to specific scene domains. Abnormalities are themselves defined in a somewhat subjective form, sometimes according to what the algorithms can detect. In some cases, different authors even define different anomalies on common data sets. Finally, experimental results can be presented on data sets of very different characteristics (e.g., traffic intersection versus subway entrance), frequently proprietary, and with widely varying levels of crowd density.

In this work, we propose an integrated solution to all these problems. We start by introducing normalcy models that *jointly account for the appearance and dynamics of complex crowd scenes*. This is done by resorting to a video representation based on dynamic textures (DTs) [11]. This representation is then used to design models of normalcy over both space and time. *Temporal normalcy* is modeled with a mixture of DTs [12] (MDT) and enables the detection of behaviors that deviate from those observed in the past. *Spatial normalcy* is measured with a discriminant saliency detector [13] based on MDTs, enabling the detection of behaviors that deviate from those of the surrounding crowd. The integration of spatial and temporal normalcy

• W. Li and N. Vasconcelos are with the Electrical and Computer Engineering Department, University of California, San Diego, 9500 Gilman Drive, La Jolla, CA 92093.
E-mail: {wel017, nvasconcelos}@ucsd.edu.

• V. Mahadevan is with Yahoo! Labs, Embassy Golf Links Business Park, Bengaluru 560071, India. E-mail: vijay.mahadevan@gmail.com.

Manuscript received 15 Apr. 2012; revised 26 Feb. 2013; accepted 14 May 2013; published online 12 June 2013.

Recommended for acceptance by G. Mori.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2012-04-0294.

Digital Object Identifier no. 10.1109/TPAMI.2013.111.

with respect to either appearance or dynamics leads to a flexible model of normalcy, applicable to the detection of anomalies of relevance to various surveillance tasks.

To address the *scale* problem, MDTs are learned at multiple spatial scales. This is done with an efficient hierarchical model, where layers of MDTs with successively larger regions of video support are learned recursively. The local measures of spatial and temporal abnormality are then integrated into a globally coherent anomaly map, by probabilistic inference. This is implemented with a conditional random field (CRF), whose single-node potentials are classifiers of local measures of spatial and temporal abnormality, collected over a range of spatial scales. They are complemented by a novel set of interaction potentials, which account for spatial and temporal context, and integrate anomaly information across the visual field.

Finally, to address the difficulties of empirical evaluation of anomaly detectors on crowded scenes, we introduce a *data set* of video from walkways in the campus of University of California, San Diego (UCSD), depicting crowds of varying densities. The data set contains 98 video sequences, and five well-defined abnormal categories. These are not “synthetic,” or “staged,” but abnormal events that occur naturally, for example, bicycle riders that cross pedestrian walkways. Ground truth is provided for abnormal events, as well as a protocol to evaluate detection performance.

The remainder of the paper is organized as follows: Section 2 reviews previous work on anomaly detection in computer vision. The problems of temporal and spatial anomaly detection in crowded scenes are discussed in Section 3. This is followed by the mathematical characterization of multiscale anomaly maps in Section 4, and the proposed CRF for integration of spatial and temporal anomalies across different spatial scales in Section 5. Finally, an extensive experimental evaluation is discussed in Section 6 and some conclusions are presented in Section 7.

2 PRIOR WORK

Recent advances in anomaly detection address event representation and globally consistent statistical inference. Contributions of the first type define features and models for the discrimination of normal and anomalous patterns. Models of normal and abnormal behavior are then learned from training data, and anomalies detected with a minimum probability of error decision rule. Although there are some exceptions [5], the distribution of abnormal patterns is usually assumed uniform, and abnormal events formulated as events of low probability under the model of normalcy.

One intuitive representation for event modeling is based on object trajectories. It is comprised of either explicitly or implicitly segmenting and tracking each object in the scene, and fitting models to the resulting object tracks [14], [15], [16], [6], [17], [18]. While capable of identifying abnormal behaviors of high-level semantics (e.g., unusual long-term trajectories), these procedures are both difficult and computationally expensive for crowded or cluttered scenes. A number of promising alternatives, which avoid processing individual objects, have been recently proposed. These include the modeling of motion patterns with histograms of pixel change [5], histograms of optical flow [19], [8], [20], or optical flow measures [3], [4], [17], [1]. Among these, [3]

models local optical flow with a mixture of probabilistic principal component analysis (PCA) models, [4] and [17] draw inspiration from classical studies of crowd behavior [21] to characterize flow with interaction features (e.g., social force model), and [1] learns the representative flow of groups by clustering optical flow-based particle trajectories.

These approaches emphasize dynamics, ignoring anomalies of object appearance and, thus, anomalous behavior without outlying motion. Optical flow, pixel change histograms, or other classical background subtraction features are also difficult to extract from crowded scenes, where the background is by definition dynamic, there are lots of clutter, and occlusions. More complete representations account for both appearance and motion. For example, [2] models temporal sequences of spatiotemporal gradients to detect anomalies in densely crowded scenes, [22] declares as abnormal spatiotemporal patches that cannot be reconstructed from previous frames, and [23] pools appearance and motion features over spatial neighborhoods, using a distance to the nearest spatially collocated feature vector among all training video clips, to quantify abnormality.

Object-based representations, based on location, blob shape, and motion [7] or optical flow magnitude, gradients, location, and scale [9], have also been proposed. Other representations include a bag-of-words over a set of manually annotated event classes [24]. Various methods have also been used to produce anomaly scores. While simple spatial filtering suffices for some applications [19], crowded scenes require more sophisticated graphical models and inference. For example, [6] and [1] adopt Gaussian mixture models (GMM) to represent trajectories of normal behavior. Cong et al. [8] and Zhao et al. [20] learn a sparse basis and define unusual events as those that can only be reconstructed with either large error or the combination of a large number of basis vectors.

Contributions of the second type address the integration of local anomaly scores, which can be noisy, into a globally consistent anomaly map. The authors of [2], [25], and [7] guarantee temporally consistent inference by modeling normal temporal sequences with hidden Markov models (HMMs). While this enforces consistency along the temporal dimension, there have also been efforts to produce spatially consistent anomaly maps. For example, latent Dirichlet allocation (LDA) has been applied to force flow features, in the model of spatial crowd interactions of [4]. On the other hand, [5] and [3] rely on Markov random fields (MRF) to enforce global spatial consistency. In the realm of sparse representations, [20] guarantees consistency of reconstruction coefficients over space and time by inclusion of smoothness terms in the underlying optimization problem. Finally, [9] models object relationships, using Bayesian networks to implement occlusion reasoning.

It should be noted that most of these methods have not been tested on the densely crowded scenes considered in this work. It is unclear that many of them could deal with the complex motion and object interactions prevalent in such scenes. Furthermore, while most methods include some mechanism to encourage spatial and temporal consistency of anomaly judgments (MRF, LDA, etc.), the underlying decision rule tends to be either predominantly temporal (e.g., trajectories, GMMs, HMMs, or sparse representations learned over time) or spatial

(e.g., interaction models) but is rarely discriminant with respect to both space and time. This makes it difficult to infer whether spatial or temporal modeling are critically important by themselves, or what benefits are gained from their joint modeling. Furthermore, the role of scale is rarely considered. These issues motivate the contributions of the following sections.

3 ANOMALY DETECTION

We start by proposing an anomaly detector that accounts for scene appearance and dynamics, spatial and temporal context, and multiple spatial scales.

3.1 Mathematical Formulation

A classical formulation of anomaly detection, which we adopt in this work, equates anomalies to outliers. A statistical model $p_{\mathbf{X}}(\mathbf{x})$ is postulated for the distribution of a measurement \mathbf{X} under *normal* conditions. Abnormalities are defined as measurements whose probability is below a threshold under this model. This is equivalent to a statistical test of hypotheses:

- \mathcal{H}_0 : \mathbf{x} is drawn from $p_{\mathbf{X}}(\mathbf{x})$;
- \mathcal{H}_1 : \mathbf{x} is drawn from an uninformative distribution $p_{\mathbf{X}}(\mathbf{x}) \propto 1$.

The minimum probability of error rule for this test is to reject the null hypothesis \mathcal{H}_0 if $p_{\mathbf{X}}(\mathbf{x}) < \nu$, where ν is the normalization constant of the uninformative distribution. As usual in the literature, we consider the problem of anomaly detection from *localized* video measurements \mathbf{x} , where \mathbf{x} is a spatiotemporal patch of small dimensions.

3.2 Spatial versus Temporal Anomalies

The normalcy model $p_{\mathbf{X}}(\mathbf{x})$ can have both a *temporal* and a *spatial* component. Temporal normalcy reflects the intuition that normal events are *recurrent* over time, i.e., previous observations establish a contextual reference for normalcy judgments. Consider a highway lane where cars move with a certain orientation and speed. Bicycles or cars heading in the opposite direction are easily identified as abnormal because they give rise to observations \mathbf{x} substantially different from those collected in the past. In this sense, temporal normalcy detection is similar to *background subtraction* [26]. A model of normal behavior is learned over time, and measurements that it cannot explain are denoted *temporal anomalies*.

Spatial normalcy reflects the intuition that some events that would not be abnormal per se are abnormal *within* a crowd. Since the crowd places physical or psychological constraints on individual behavior, behaviors feasible in isolation can have low probability in a crowd context. For example, while there is nothing abnormal about an ambulance that rides at 50 mph in a stretch of highway, the same observation within a highly *congested* highway is abnormal. Note that the only indication of abnormality is the *difference* between the crowd and the object at the time of the observation, not that the ambulance moves at 50 mph. Since the detection of such abnormalities is mostly based on spatial context, they are denoted *spatial anomalies*. Their detection does not depend on memory. Instead, it is based on a continuously evolving, instantaneously adaptive,

definition of normalcy. In this sense, the detection of spatial anomalies can be equated to *saliency detection* [27].

3.3 Roles of Crowds and Scale

Most available background subtraction and saliency detection solutions are not applicable to crowded scenes, where backgrounds can be highly dynamic. In this case, it is not sufficient to detect variations of image intensity, or even optical flow, to detect anomalous events. Instead, normalcy models must rely on sophisticated *joint representations of appearance and dynamics*. In fact, even such models can be ineffective. Since crowds frequently contain distinct sub-entities, for example, vehicles or groups of people moving in different directions, anomaly detection requires modeling *multiple video components of different appearance and dynamics*. A model that has been shown successful in this context is the mixture of DTs [12]. This is the representation adopted in this work.

Another challenging aspect of anomaly detection within crowds is scale. Spatial anomalies are usually detected at the scale of the smallest scene entities, typically people. However, a normal event at this scale may be anomalous at a larger scale, and vice versa. For example, while a child that rides a bicycle appears normal within a group of bicycle riding children, the group is itself anomalous in a crowded pedestrian sidewalk. Local anomaly detectors, with small regions of interest, cannot detect such anomalies. To address this, we represent crowded scenes with a hierarchy of MDTs that cover successively larger regions. This is done with a computationally efficient hierarchical model, where MDT layers are estimated recursively.

A similar challenge holds for temporal anomalies. While their detection is usually based on a small number of video frames, certain anomalies can only be detected over long time spans. For example, while it is normal for two pedestrian trajectories to converge or diverge at any point in time, a cyclical convergence and divergence is probably abnormal. Anomaly detection across time scales is, however, more complex than across spatial scales, due to constraints of instantaneous detection and implementation complexity. Since video has to be buffered before anomalies can be detected, large temporal windows imply long detection delays and storage of many video frames. Due to this, we do not consider multiple temporal scales in this work. A single scale is chosen, using acceptable values of delay and storage complexity, and used throughout our experiments. Note that, like their spatial counterparts, temporal anomaly maps are computed at multiple spatial scales. Hence, in what follows, the term “scale” refers to the spatial support of anomaly detection, for *both* spatial and temporal anomalies.

4 NORMALCY AND ANOMALY MODELING

In this section, we review the MDT model, discuss the design of temporal and spatial models of normalcy, and formulate the computation of anomaly maps.

4.1 Mixture of Dynamic Textures

The MDT models a sequence of τ video frames $\mathbf{x}_{1:\tau} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\tau]$ as a sample from one of K dynamic textures [11]:

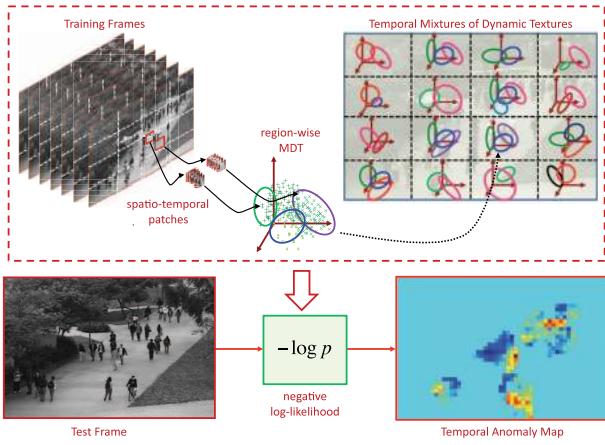


Fig. 1. Temporal anomaly detection. An MDT is learned per scene subregion, at training time. A temporal anomaly map is produced by measuring the negative log probability of each video patch under the MDT of the corresponding region.

$$p(\mathbf{x}_{1:\tau}) = \sum_{i=1}^K \pi_i p(\mathbf{x}_{1:\tau} | z = i). \quad (1)$$

The mixture components $p(\mathbf{x}_{1:\tau} | z = i)$ are linear dynamic systems (LDS) defined by

$$\begin{cases} \mathbf{s}_{t+1} = A_z \mathbf{s}_t + \mathbf{n}_t, & (2a) \\ \mathbf{x}_t = C_z \mathbf{s}_t + \mathbf{m}_t, & (2b) \end{cases}$$

where Z is a multinomial random variable of parameters $\boldsymbol{\pi}$ ($\pi_i \geq 0, \sum_i \pi_i = 1$), which indexes the mixture component from which \mathbf{x}_t is drawn. \mathbf{s}_t is a hidden state variable that encodes scene dynamics, and \mathbf{x}_t the vector of pixels in video frame t . A_z, C_z are the transition and observation matrices of component z , whose initial condition is $\mathbf{s}_1 \sim \mathcal{N}(\boldsymbol{\mu}_z, S_z)$, and noise processes are defined by $\mathbf{n}_t \sim \mathcal{N}(\mathbf{0}, Q_z)$ and $\mathbf{m}_t \sim \mathcal{N}(\mathbf{0}, R_z)$. The model parameters are learned by maximum-likelihood estimation (MLE) from a collection of video patches, with the expectation-maximization (EM) algorithm of [12], which is reviewed in Appendix A.1, which is available in the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TPAMI.2013.111>.

4.2 Temporal Anomaly Detection

Temporal anomaly detection is inspired by the popular background subtraction method of [26]. This uses a GMM per image location to model the distribution of image intensities. Observations of low probability under these GMMs are declared foreground. For anomaly detection in crowds, the GMM is replaced by an MDT, and the pixel grid replaced by one of preset displacement. Grid locations define the center of video cells, from which video patches are extracted. The patches extracted from a subregion (group of cells) are used to learn an MDT, during a training phase, as illustrated in Fig. 1. After this phase, subregion patches of low probability under the associated MDT are considered anomalies. Given patch $\mathbf{x}_{1:\tau}$, the distribution of the hidden state sequence \mathbf{s}_1^{τ} under the i th DT component, $p_{S|\mathbf{X}}(\mathbf{s}_{1:\tau} | \mathbf{x}_{1:\tau}, z = i)$, is estimated with a Kalman filter and smoother [28], [29], as discussed in Appendix A.2, available in the online supplemental material. The value of the temporal anomaly

map at location l is the negative-log probability of the most-likely state sequence for the patch at l :

$$\mathcal{T}(l) = -\log \left[\sum_{i=1}^K \pi_i p(\mathbf{s}_{1:\tau}^{\{i\}}(l) | z = i) \right], \quad (3)$$

where $\mathbf{s}_{1:\tau}^{\{i\}}(l) = \operatorname{argmax}_{\mathbf{s}_{1:\tau}} p(\mathbf{s}_{1:\tau} | \mathbf{x}^{\tau}(l), z = i)$. We note that this generalizes the mixture of PCA models of optical flow [3]. The matrix C_z of (2b) is a PCA basis for patches drawn from mixture component z , but the PCA decomposition reports to *patch appearance*, not optical flow. Patch dynamics are captured by the hidden state sequence $\mathbf{s}_{1:\tau}$, which is a trajectory in PCA space. Hence, unlike mixtures of optical flow, the representation is temporally *smooth*. The joint representation of appearance and dynamics makes the MDT a better representation for crowd video than the mixture of PCA.

4.3 Spatial Anomaly Detection

Spatial anomaly detection is inspired by previous work in saliency detection [27], [13]. Saliency is defined in a center-surround manner. Given a set of features, salient locations are those of substantial feature contrast with their immediate surround. Spatial anomalies are then defined as locations whose saliency is above some threshold. In this work, we rely on the discriminant saliency criterion of [13].

4.3.1 Discriminant Saliency

Discriminant saliency formulates the saliency problem as a hypothesis test between two classes: a class of *salient stimuli*, and a *background* class of stimuli that are not salient. Two windows are defined at each scene location l : a *center window* \mathcal{W}_l^1 , with label $\mathcal{C}(l) = 1$, containing the location, and a *surrounding annular window* \mathcal{W}_l^0 , with label $\mathcal{C}(l) = 0$, containing *background*. A set of feature responses \mathbf{X} are computed for each of the windows \mathcal{W}_l^c , $c \in \{0, 1\}$ and $\mathcal{S}(l)$, the saliency of location l , defined as the extent to which they discriminate between the two classes. This is quantified by the mutual information (MI) between feature responses and class label [13]:

$$\mathcal{S}(l) = \sum_{c=0}^1 \{p_{\mathcal{C}(l)}(c) \text{KL}[p_{\mathbf{X}|\mathcal{C}(l)}(\mathbf{x}|c) \| p_{\mathbf{X}}(\mathbf{x})]\}, \quad (4)$$

where $p_{\mathbf{X}|\mathcal{C}(l)}(\mathbf{x}|c)$ are class-conditional densities and $\text{KL}(p||q) = \int_{\mathcal{X}} p_{\mathbf{X}}(\mathbf{x}) \log \frac{p_{\mathbf{X}}(\mathbf{x})}{q_{\mathbf{X}}(\mathbf{x})} d\mathbf{x}$ the Kullback-Leibler (KL) divergence between $p_{\mathbf{X}}(\mathbf{x})$ and $q_{\mathbf{X}}(\mathbf{x})$ [30].

Locations of maximal saliency are those where the discrimination between center and surround can be made with highest confidence, i.e., where (4) is maximal. The discriminant saliency principle can be applied to many features [31]. When \mathbf{X} consists of optical flow, it generalizes the force flow model of [4], where saliency is defined as the difference between the optical flow at l and the average flow in its neighborhood (see [4, (8)]). This is a simplified form of discriminant saliency, which replaces the MI of (4) by a difference to the mean background response.

4.3.2 Center-Surround Saliency with MDTs

Optical flow methods provide a coarse representation of dynamics and ignore appearance. For background subtraction, this problem has been addressed with the combination of DTs and discriminant saliency [32]. While using a more

powerful representation than force flow, this method learns a single DT from both center and surround windows. This assumes a homogeneity of appearance and dynamics within the two windows that do not hold for crowds, where foregrounds and backgrounds can be quite diverse.

In this work, we adopt the MDT as the probability distribution $p_{\mathbf{X}|c(l)}(\mathbf{x}_{1:\tau}|c)$ from which spatiotemporal patches $\mathbf{x}_{1:\tau}^c$ are drawn. We note that under assumptions of Gaussian initial conditions and noise, patches $\mathbf{x}_{1:\tau}$ drawn from a DT have a Gaussian probability distribution [33],

$$\mathbf{x}_{1:\tau} \sim \mathcal{N}(\boldsymbol{\gamma}, \Phi), \quad (5)$$

whose parameters follow from those of the LDS (2). When the class-conditional distributions of the center and surround classes, $c \in \{0, 1\}$, at location l are mixtures of K_c DTs, it follows that

$$\begin{aligned} p_{\mathbf{X}|c(l)}(\mathbf{x}_{1:\tau}|c) &= \sum_{i=1}^{K_c} \pi_i^c \mathcal{N}(\mathbf{x}_{1:\tau}; \boldsymbol{\gamma}_i^c, \Phi_i^c) \\ &= \sum_{i=1}^{K_c} \pi_i^c p_{\mathbf{X}|c(l)}^i(\mathbf{x}_{1:\tau}|c), \end{aligned} \quad (6)$$

for $c \in \{0, 1\}$. The marginal distribution is then

$$\begin{aligned} p_{\mathbf{X}}(\mathbf{x}_{1:\tau}) &= \sum_{c=0}^1 [p_{c(l)}(c) p_{\mathbf{X}|c(l)}(\mathbf{x}_{1:\tau}|c)] \\ &= \sum_{c=0}^1 \left[p_{c(l)}(c) \sum_{i=1}^{K_c} \pi_i^c \mathcal{N}(\mathbf{x}_{1:\tau}; \boldsymbol{\gamma}_i^c, \Phi_i^c) \right] \\ &= \sum_{i=1}^{K_0+K_1} \omega_i \mathcal{N}(\mathbf{x}_{1:\tau}; \boldsymbol{\gamma}_i, \Phi_i) \\ &= \sum_{i=1}^{K_0+K_1} \omega_i p_{\mathbf{X}}^i(\mathbf{x}_{1:\tau}), \end{aligned} \quad (7)$$

and the saliency measure of (4) requires the KL divergence between (6) and (7). This is problematic because there is no closed form solution for the KL divergence between two MDTs. However, because the MDT components are Gaussian, it is possible to rely on popular approximations to the KL divergence between Gaussian mixtures. We adopt the variational approximation of [34]:

$$\begin{aligned} \text{KL}(p_{\mathbf{X}|c} \| p_{\mathbf{X}}) \\ \approx \sum_i \left\{ \pi_i^c \log \frac{\sum_j^{K_c} \pi_j^c \exp(-\text{KL}(p_{\mathbf{X}|c}^i \| p_{\mathbf{X}|c}^j))}{\sum_j^{K_0+K_1} \omega_j \exp(-\text{KL}(p_{\mathbf{X}|c}^i \| p_{\mathbf{X}}^j))} \right\}. \end{aligned} \quad (8)$$

Each term of (8) contains a KL divergence between DTs, which can be computed in closed form [35]. For example, for the terms in the denominator

$$\begin{aligned} \text{KL}(p_{\mathbf{X}|c}^i \| p_{\mathbf{X}}^j) \\ = \frac{1}{2} \left[\log \frac{|\Phi_j|}{|\Phi_i^c|} + \text{Tr}(\Phi_j^{-1} \Phi_i^c) + \|\boldsymbol{\gamma}_i^c - \boldsymbol{\gamma}_j\|_{\Phi_j}^2 - m\tau \right], \end{aligned} \quad (9)$$

where m is the number of pixels per frame, and $\|\mathbf{z}\|_{\Phi} = \mathbf{z}^T \Phi^{-1} \mathbf{z}$. Numerator terms are computed similarly. All computations can be performed recursively [35].

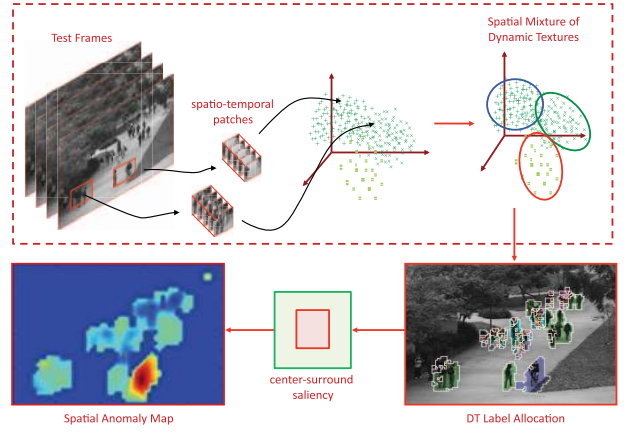


Fig. 2. Spatial anomaly detection using center-surround saliency with MDT models.

4.3.3 Spatial Anomaly Map

The spatial anomaly map is a map of the saliency $S(l)$ at locations l . Given a location, this requires 1) learning MDTs from center and surround windows, and 2) computing a weighted average of these mixtures to obtain (7). Since learning MDTs per location is computationally prohibitive, we resort to the following approximation. A dense collection of overlapping spatiotemporal patches is first extracted from $\mathcal{V}(t)$, a 3D video volume temporally centered at the current frame. A single MDT with K^g mixture components, denoted $\{\boldsymbol{\gamma}_i^g, \Phi_i^g\}_{i=1}^{K^g}$, is learned from this patch collection. Each patch is then assigned to the mixture component of largest posterior probability. This segments the volume into superpixels, as shown in Fig. 2.

At location l , the MDTs of (6) and (7) are derived from the global mixture model. The DT components are assumed equal to those of the latter and only the mixing proportions are recomputed, using the ratio of pixels assigned to each component in the respective windows:

$$p_{\mathbf{X}|c(l)}(\mathbf{x}_{1:\tau}|c) = \sum_{i=1}^{K^g} \frac{\sum_{l \in \mathcal{W}_i^c} \mathcal{M}_{il}}{\sum_{l \in \mathcal{W}_i^c} 1} \mathcal{N}(\mathbf{x}_{1:\tau}; \boldsymbol{\gamma}_i^g, \Phi_i^g), \quad (10)$$

for $c \in \{0, 1\}$. $\mathcal{M}_{il} = 1$ if l is assigned to mixture component i and 0 otherwise. The prior probabilities for center and surround, $p_c(c)$, are proportional to the ratio of volumes of center and surround windows. $S(l)$ is computed with (4), using (8) and (9). Note that the KL divergence terms in (8) only require the computation of $\binom{K^g}{2}$ KL divergences between the K^g mixture components, and these are computed only once per frame because all mixture components are shared (i.e., the terms $\exp(-\text{KL}(p\|q))$ in (8) are fixed per frame). This procedure is repeated for every frame in the test video, as illustrated in Fig. 2.

4.4 Multiscale Anomaly Maps

To account for anomalies at multiple spatial scales, we rely on a hierarchical mixture of dynamic textures (H-MDT). This is a model with various MDT layers, learned from regions of different spatial support. At the finest scale, a video sequence is divided into n_L subregions (e.g., 5×8

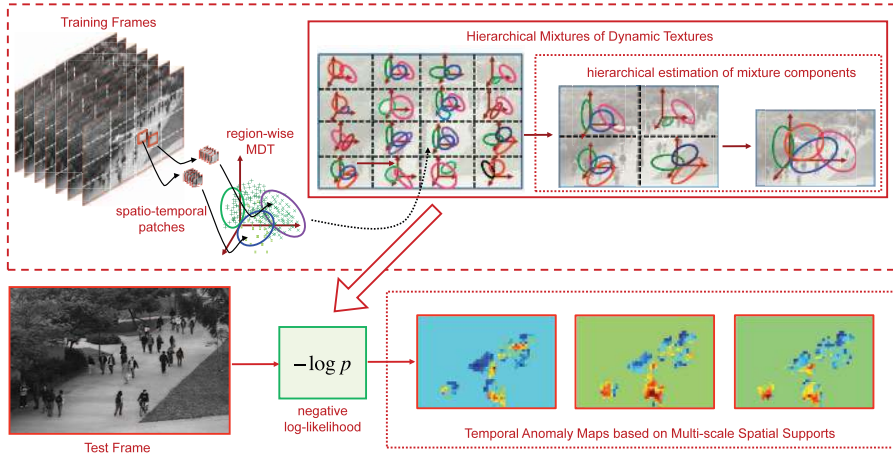


Fig. 3. Computation of temporal anomaly maps with multiscale spatial supports using the H-MDT. MDTs of increasingly larger spatial support are estimated recursively, with the H-EM algorithm. Their application to a query video produces temporal anomaly maps based on supports of various spatial scales.

subregions). n_L MDT models $\{\mathbf{M}_i\}_{i=1}^{n_L}$ are then learned from patches extracted from each of the subregions. At the coarsest scale, the whole visual field is represented with a global MDT. This results in a hierarchy of MDT models $\{\{\mathbf{M}_i^1\}_{i=1}^{n_1}, \dots, \mathbf{M}_1^L\}$, where \mathbf{M}_j^s , the j th model at scale s , is learned from subregion \mathcal{R}_j^s . The hierarchy of support windows $\{\{\mathcal{R}_i^1\}_{i=1}^{n_1}, \dots, \mathcal{R}^L\}$ resembles the spatial pyramid structure of [36]. H-MDT models can be learned efficiently with the hierarchical expectation-maximization (H-EM) algorithm of [37]. Rather than collecting patches anew from larger regions, it estimates the models at a given layer directly from the parameters of the MDT models at the layer of immediately higher resolution.

For anomaly detection, each model is applied to the corresponding window. This produces L anomaly maps, $\{\mathcal{T}^1, \dots, \mathcal{T}^L\}$, as illustrated in Fig. 3. A hierarchy of spatial anomaly maps, $\{\mathcal{S}^1, \dots, \mathcal{S}^L\}$ is also computed. For all s , the computation of \mathcal{S}^s relies on a global mixture model \mathbf{M} . The mixing proportions of (10) are computed using surround windows of size identical to $\{\mathcal{R}_i^s\}$ and center windows of constant size, as summarized in Algorithm 1 (see Appendix B for all algorithms, available in the online supplemental material).

5 GLOBALLY CONSISTENT ANOMALY MAPS

In this section, we introduce a layer of statistical inference to fuse anomaly information across time, space, and scale in a globally consistent manner.

5.1 Discriminative Model

The anomaly maps of the previous section span space, time, and spatial scale. Being derived from local measurements, they can be noisy. A principled framework is required to 1) integrate anomaly scores from the individual maps, 2) eliminate noise, and 3) guarantee spatiotemporal consistency of anomaly judgments throughout the visual field. For this, we rely on a conditional random field [38] inspired by the discriminative random field (DRF) of [39]. An anomaly label $y_i \in \{-1, 1\}$ is defined at each location i in a set S of observation sites. Given a video clip \mathbf{x} , the

conditional likelihood of observing a configuration of anomaly labels $\mathbf{y} = \{y_i | i \in S\}$ is

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z} \exp \left\{ \sum_{i \in S} \mathcal{A}(y_i, \mathbf{x}) + \sum_{i \in S} \left[\frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathcal{I}(y_i, y_j, \mathbf{x}, i, j) \right] \right\}, \quad (11)$$

where Z is a partition function and \mathcal{N}_i the neighborhood of site i . The single-site and interaction potentials of (11),

$$\mathcal{A}(y_i, \mathbf{x}) = \log \sigma(y_i \mathbf{w}^T \mathbf{f}_i), \quad (12)$$

where $\sigma(x) = (1 + e^{-x})^{-1}$ is the sigmoid function, and

$$\mathcal{I}(y_i, y_j, \mathbf{x}, i, j) = y_i y_j \cdot \mathbf{v}^T \boldsymbol{\mu}(\mathbf{f}_i, \mathbf{f}_j, i, j) \quad (13)$$

are based on a feature vector \mathbf{f}_i that concatenates the spatial and temporal anomaly scores of site i at the L spatial scales, plus a bias term (set to 1):

$$\mathbf{f}_i = [1, \mathcal{T}^1(i), \dots, \mathcal{T}^L(i), \mathcal{S}^1(i), \dots, \mathcal{S}^L(i)]^T. \quad (14)$$

\mathbf{w} , \mathbf{v} are parameter vectors and $\boldsymbol{\mu}$ a compound feature:

$$\boldsymbol{\mu}(\mathbf{f}_i, \mathbf{f}_j, i, j) = e^{-\alpha|i-j|} \exp(-\mathbf{h}_{i,j}), \quad (15)$$

where $|i - j|$ is the euclidean distance between sites i, j , and $\exp(-\mathbf{h}_{i,j})$ the entry-wise exponential of $-\mathbf{h}_{i,j}$. The vector $\mathbf{h}_{i,j}$ contains the diagonal entries of $(\mathbf{f}_i - \mathbf{f}_j)(\mathbf{f}_i - \mathbf{f}_j)^T$.

The single-site potential of (12) reflects the anomaly belief at site i . Using it alone, i.e., without (13), (11) is a logistic regression model. In this case, the detection of each anomaly is based on information from site i exclusively. The addition of the interaction potential of (13) enables the model to take into account information from site i 's neighborhood \mathcal{N}_i . This smooths the single-site prediction, encouraging consistency of neighboring labels. The interaction potential can be interpreted as a classifier that predicts whether two neighboring sites have the same label. Note that because \mathbf{f} contains anomaly scores at different spatial scales, $\mathbf{h}_{i,j}$ (or $\boldsymbol{\mu}_{i,j}$) accounts for the

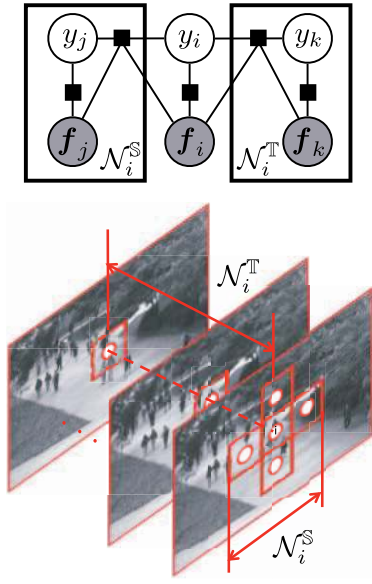


Fig. 4. CRF filter. Top: Graphical model. Bottom: Spatial and temporal neighborhoods.

similarity between the two observations in anomaly spaces of different scale (i.e., under different spatial normalcy contexts). The interaction potentials adaptively modulate the intensity of intersite smoothing according to these similarity measures (and how they are weighted by \mathbf{v}). The parameters \mathbf{w} and \mathbf{v} encode the relative importance of different features.

5.2 Online CRF Filter

The model of (11) requires inference over the entire video sequence. This is not suitable for online applications. An online version can be implemented by conditioning the anomaly label $\mathbf{y}^{(\tau)}$ at time τ on 1) observations for $t \leq \tau$, and 2) anomaly labels for $t < \tau$, leading to

$$\begin{aligned}
 & P(\mathbf{y}^{(\tau)} | \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}; \Theta) \\
 &= \frac{1}{\mathcal{Z}} \exp \left\{ \sum_{i \in S^\tau} \left[\mathcal{A}(y_i^{(\tau)}, \mathbf{x}^{(\tau)}) \right. \right. \\
 & \quad \left. \left. + \frac{1}{|\mathcal{N}_i^S|} \sum_{j \in \mathcal{N}_i^S} \mathcal{I}_S(y_i^{(\tau)}, y_j, \mathbf{x}^{(\tau)}, i, j) \right. \right. \\
 & \quad \left. \left. + \frac{1}{|\mathcal{N}_i^T|} \sum_{k \in \mathcal{N}_i^T} \mathcal{I}_T(y_i^{(\tau)}, y_k, \mathbf{x}, i, k) \right] \right\}, \quad (16)
 \end{aligned}$$

where S^τ is the set of observations at time τ (pixels of the current frame). Two neighborhoods are defined per location i : spatial \mathcal{N}_i^S ($\mathcal{N}_i^S \subseteq S^\tau$) and temporal \mathcal{N}_i^T ($\mathcal{N}_i^T \subseteq \{S^t\}_{t=1}^{\tau-1}$). The graphical model is shown at the top of Fig. 4, and these neighborhoods at the bottom. The parameters $\Theta = \{\mathbf{w}, \mathbf{v}_T, \mathbf{v}_S, \alpha_T, \alpha_S\}$ are estimated during training.

5.2.1 Learning

Both (11) and (16) can be learned with standard optimization techniques, such as gradient descent or the *Broyden-Fletcher-Goldfarb-Shanno* (BFGS) method. To improve generalization, the model is regularized with a Gaussian prior of standard

deviation ϵ , for all parameters. Given N independent training samples $\{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}\}_{n=1}^N$, the gradients of the regularized log-likelihood with respect to \mathbf{w} , \mathbf{v} , and α are

$$\begin{aligned}
 & \frac{\partial}{\partial \mathbf{w}} \log p(\{\mathbf{y}^{(n)}\}_{n=1}^N | \{\mathbf{x}^{(n)}\}_{n=1}^N) \\
 &= \sum_{n=1}^N \left\{ \sum_{i \in S} \sigma(-y_i^{(n)} \mathbf{w}^T \mathbf{f}_i^{(n)}) y_i^{(n)} \mathbf{f}_i^{(n)} \right. \\
 & \quad \left. - \mathbb{E} \left[\sum_{i \in S} \sigma(-y_i \mathbf{w}^T \mathbf{f}_i^{(n)}) y_i \mathbf{f}_i^{(n)} \right] \right\} - \frac{1}{\epsilon_{\mathbf{w}}^2} \mathbf{w}, \quad (17)
 \end{aligned}$$

$$\begin{aligned}
 & \frac{\partial}{\partial \mathbf{v}} \log p(\{\mathbf{y}^{(n)}\}_{n=1}^N | \{\mathbf{x}^{(n)}\}_{n=1}^N) \\
 &= \sum_{n=1}^N \left\{ \sum_{i \in S} \left[\frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} (e^{-\alpha|i-j|} y_i^{(n)} y_j^{(n)} \exp(-\mathbf{h}_{i,j}^{(n)})) \right] \right. \\
 & \quad \left. - \mathbb{E} \left[\sum_{i \in S} \left(\frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} (e^{-\alpha|i-j|} y_i y_j \exp(-\mathbf{h}_{i,j}^{(n)})) \right) \right] \right\} - \frac{1}{\epsilon_{\mathbf{v}}^2} \mathbf{v}, \quad (18)
 \end{aligned}$$

and

$$\begin{aligned}
 & \frac{\partial}{\partial \alpha} \log p(\{\mathbf{y}^{(n)}\}_{n=1}^N | \{\mathbf{x}^{(n)}\}_{n=1}^N) \\
 &= \sum_{n=1}^N \left\{ \sum_{i \in S} \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} (-\mathcal{I}(y_i^{(n)}, y_j^{(n)}, \mathbf{x}^{(n)}, i, j) |i-j|) \right. \\
 & \quad \left. + \mathbb{E} \left[\sum_{i \in S} \left(\frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} (\mathcal{I}(y_i, y_j, \mathbf{x}^{(n)}, i, j) |i-j|) \right) \right] \right\} - \frac{1}{\epsilon_{\alpha}^2} \alpha, \quad (19)
 \end{aligned}$$

where the expectation is evaluated with distribution $p(\mathbf{Y} | \mathbf{X}; \Theta)$. The conditional expectations of (17)-(19) require evaluation of the partition function \mathcal{Z} , a problem known to be NP-hard. As is common in the literature, this difficulty is avoided by estimating expectations through sampling. Although sampling methods such as *Markov chain Monte Carlo* (MCMC) can converge to the true distribution, this usually requires many iterations. Since the procedure must be repeated per gradient ascent step, these methods are impractical. On the other hand, approximations such as contrastive divergence minimization (which runs MCMC a limited number of times with specific starting points) have been shown to be successful for vision applications [40], [41]. We adopt these approximations for CRF learning.

This leverages the fact that, denoting any of the parameters $\mathbf{w}, \mathbf{v}_T, \mathbf{v}_S, \alpha_T, \alpha_S$ by θ , the partial gradients of (17)-(19) are

$$\begin{aligned}
 & \frac{\partial}{\partial \theta} \log p(\{\mathbf{y}^{(n)}\}_{n=1}^N | \{\mathbf{x}^{(n)}\}_{n=1}^N; \Theta) \\
 &= \sum_{n=1}^N \left\{ F_{\partial \theta}(\mathbf{y}^{(n)}, \mathbf{x}^{(n)}) - \mathbb{E}_{(\mathbf{Y} | \mathbf{X}; \Theta)} [F_{\partial \theta}(\mathbf{y}, \mathbf{x}^{(n)})] \right\} - \frac{1}{\epsilon_{\theta}^2} \theta, \quad (20)
 \end{aligned}$$

where $F_{\partial \theta}(\mathbf{y}, \mathbf{x})$ is the sum of the terms in the summations of (17), (18), or (19) that depend on θ . Contrastive divergence approximates the intractable conditional

expectation $\mathbb{E}_{(\mathbf{Y}|\mathbf{X},\Theta)}[F_{\partial\theta}(\mathbf{y}, \mathbf{x}^{(n)})]$ by $F_{\partial\theta}(\hat{\mathbf{y}}, \mathbf{x}^{(n)})$, where $\hat{\mathbf{y}}$ is the ‘‘evil twin’’ of the ground-truth label field $\mathbf{y}^{(n)}$ [41]. $\hat{\mathbf{y}}$ is drawn by MCMC, using the inference procedure discussed in Section 5.2.2, the current parameter estimates, and the ground-truth labels $\mathbf{y}^{(n)}$ as a starting point.

Given the estimate of the partial gradients, the gradient ascent rule for parameter updates reduces to

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \left[\sum_{n=1}^N (F_{\partial\theta}(\mathbf{y}^{(n)}, \mathbf{x}^{(n)}) - F_{\partial\theta}(\hat{\mathbf{y}}^{(n)}, \mathbf{x}^{(n)})) - \frac{1}{\epsilon_{\boldsymbol{\theta}}^2} \boldsymbol{\theta} \right], \quad (21)$$

where η is a learning rate. In our implementation, this rule is initialized with $\mathbf{v}_{\mathbb{T}} = \mathbf{v}_{\mathbb{S}} = \mathbf{1}$ and $\alpha_{\mathbb{T}} = \alpha_{\mathbb{S}} = \mathbf{0}$. The initial value of \mathbf{w} is learned, assuming a logistic regression model ($\mathbf{v}_{\mathbb{T}} = \mathbf{v}_{\mathbb{S}} = \mathbf{0}$ in (16)), with the procedure of [43].

5.2.2 Inference

The inference problem is to determine the most likely anomaly prediction \mathbf{y}^* for a query frame $\mathbf{x}^{(\tau)}$, given previous predictions $\{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}$, and observations $\{\mathbf{x}^{(t)}\}_{t=1}^{\tau}$:

$$\begin{aligned} \mathbf{y}^* &= \underset{\mathbf{y}}{\operatorname{argmax}} \log p(\mathbf{y} | \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}; \Theta) \\ &= \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i \in \mathbb{S}} \left[\mathcal{A}(y_i, \mathbf{x}^{(\tau)}) + \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathcal{I}(y_i, y_j, \mathbf{x}, i, j) \right]. \end{aligned} \quad (22)$$

Again, exact inference is intractable. We rely on Gibbs sampling to approximate the optimal prediction. This consists of drawing labels from the conditional distribution:

$$\begin{aligned} p(y_i | \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \mathbf{y}_{-i}; \Theta) \\ &= \frac{p(y_i, \mathbf{y}_{-i} | \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}; \Theta)}{p(\mathbf{y}_{-i} | \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}; \Theta)} \\ &= \frac{1}{\mathcal{Z}_{-i}} \exp [F_i(\{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \mathbf{y}_{-i}, y_i; \Theta)], \end{aligned} \quad (23)$$

where $F_i(\{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \mathbf{y}_{-i}, y_i; \Theta)$ is the sum of potential functions that depend on site i (i.e., its ‘‘Markov blanket’’):

$$\begin{aligned} F_i(\{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \mathbf{y}_{-i}, y_i; \Theta) \\ &= \mathcal{A}(y_i, \mathbf{x}^{(\tau)}) + \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathcal{I}(y_i, y_j, \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, i, j) \\ &\quad + \sum_{j \in \mathcal{N}_j} \frac{1}{|\mathcal{N}_j|} \mathcal{I}(y_j, y_i, \{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, j, i), \end{aligned} \quad (24)$$

and \mathcal{Z}_{-i} the corresponding partition function:

$$\mathcal{Z}_{-i} = \sum_{y_i} \exp [F_i(\{\mathbf{x}^{(t)}\}_{t=1}^{\tau}, \{\mathbf{y}^{(t)}\}_{t=1}^{\tau-1}, \mathbf{y}_{-i}, y_i; \Theta)]. \quad (25)$$

The procedure is detailed in Algorithms 2 and 3, available in the online supplemental material, where we present the online CRF filter used to estimate the label field. During learning, the filter is initialized with the ground-truth labels ($\mathbf{y}_0 = \mathbf{y}^{(\tau)}$). During testing, this initialization relies on the predictions of the single-site classifiers ($\mathbf{v}_{\mathbb{T}} = \mathbf{v}_{\mathbb{S}} = \mathbf{0}$). In our

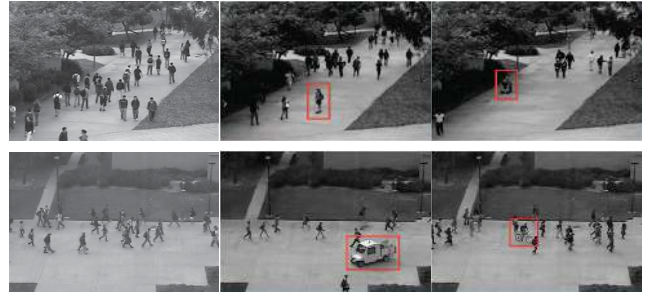


Fig. 5. Exemplar normal/abnormal frames in Ped1 (top) and Ped2 (bottom). Anomalies (red boxes) include bikes, skaters, carts, and wheelchairs.

implementation, the filter is run for $N_s = 10$ iterations. Again, the complete anomaly detection procedure is summarized in Algorithm 4, available in the online supplemental material.

6 EXPERIMENTS

In this section, we introduce a new data set and an experimental protocol for evaluation of anomaly detection in crowded environments and use it to evaluate the proposed anomaly detector.

6.1 UCSD Pedestrian Anomaly Data Set

In the literature, anomaly detection has frequently been evaluated by visual inspection [19], [7], [3], or with coarse ground truth, for example, frame-level annotation of abnormal events [4], [1]. This does not completely address the anomaly detection problem, where it is usually desired to *localize* anomalies in both *space and time*. To enable this, we introduce a data set¹ of crowd scenes with precisely localized anomalies and metrics for the evaluation of their detection. The data set consists of video clips recorded with a stationary camera mounted at an elevation, overlooking pedestrian walkways on the UCSD campus. The crowd density in the walkways is variable, ranging from sparse to very crowded. In the normal setting, the video contains only pedestrians. Abnormal events are due to either 1) the circulation of nonpedestrian entities in the walkways, or 2) anomalous pedestrian motion patterns. Commonly occurring anomalies include bikers, skaters, small carts, and people walking across a walkway or in the surrounding grass. A few instances of wheelchairs are also recorded. All abnormalities occur naturally, i.e., they were not staged or synthesized for data set collection.

The data set is organized into two subsets, corresponding to the two scenes of Fig. 5. The first, denoted ‘‘Ped1,’’ contains clips of 158×238 pixels, which depict groups of people walking toward and away from the camera, and some amount of perspective distortion. The second, denoted ‘‘Ped2,’’ has spatial resolution of 240×360 pixels and depicts a scene where most pedestrians move horizontally. The video footage of each scene is sliced into clips of 120-200 frames. A number of these (34 in Ped1 and 16 in Ped2) are to be used as training set for the condition of

1. Available from <http://www.svcl.ucsd.edu/projects/anomaly/data-set.html>.

TABLE 1
Composition of UCSD Anomaly Data Set

Scene	Nor.	Abnormal ^a					Total ^b
		Bike	Skater	Cart	Walk Across	Other	
Ped1	34	19/28	13/13	6/6	3/4	3/3	36/54
Ped2	16	11/19	3/3	1/1	0/0	0/0	12/23

^a number of clips/number of anomaly instances. ^b some clips contain more than one type of anomaly.

normalcy. The test set contains clips (36 for Ped1 and 12 for Ped2) with both normal (around 5,500) and abnormal (around 3,400) frames. The abnormalities of each set are summarized in Table 1.

Frame-level ground-truth annotation, indicating whether anomalies occur within each frame, and manually collected pixel-level binary anomaly masks, which identify the pixels containing anomalies, are available per test clip. We note that this includes ground truth on Ped1 contributed by Antić and Ommer [9], and supersedes the ground truth available on an earlier version of this work [43]. We denote the current ground truth by “full annotation” and the previous one by “partial annotation.” Unless otherwise noted, the results of the subsequent sections correspond to the full annotation.

6.2 Evaluation Methodology

Two criteria are used to evaluate anomaly detection accuracy: a *frame-level criterion* and a *pixel-level criterion*. Both are based on true-positive rates (TPR) and false-positive rates (FPRs), denoting “an anomalous event” as “positive” and “the absence of anomalous events” as “negative.” A frame containing anomalies is denoted a positive, otherwise a negative. The true and false positives under the two criteria are:

- *Frame-level criterion.* An algorithm predicts which frames contain anomalous events. This is compared to the clip’s frame-level ground-truth anomaly annotations to determine the number of true- and false-positive frames.
- *Pixel-level criterion.* An algorithm predicts which pixels are related to anomalous events. This is compared to the pixel-level ground-truth anomaly annotation to determine the number of true-positive and false-positive frames. A frame is a true positive if 1) it is positive and 2) at least 40 percent of its anomalous pixels are identified; a frame is a false positive if it is negative and any of its pixels are predicated as anomalous.

The two measures are combined into a receiver operating characteristic (ROC) curve of TPR versus FPR:

$$\text{TPR} = \frac{\# \text{ of true-positive frame}}{\# \text{ of positive frame}},$$

$$\text{FPR} = \frac{\# \text{ of false-positive frame}}{\# \text{ of negative frame}}.$$

Performance is also summarized by the *equal error rate* (EER), the ratio of misclassified frames at which $\text{FPR} = 1 - \text{TPR}$, for the frame-level criterion, or *rate of detection* (RD), i.e., $1 - \text{EER}$, for the pixel-level criterion.

Note that, although widely used in the literature, the frame-level criterion only measures temporal localization accuracy. This enables errors due to “lucky co-occurrences” of prediction errors and true abnormalities. For example, it assigns a perfect score to an algorithm that identifies a single anomaly at a random location of a frame with anomalies. The pixel-level criterion is much stricter and more rigorous. By evaluating both the temporal and spatial accuracy of the anomaly predictions, it rules out these “lucky co-occurrences.” We believe that the pixel-level criterion should be the predominant criterion for evaluation of anomaly detection algorithms.

6.3 Experimental Setup

Unless otherwise noted, observation sites are a video sublattice with spatial interval of four pixels and temporal interval of five frames. Temporal anomaly maps rely on patches of $13 \times 13 \times 15$ pixels. The temporal extent of 15 frames provides a reasonable compromise between the ability to detect anomalies and the delay (1.5 s) and storage (15 video frames) required for anomaly detection. To minimize computation, patches of variance smaller than 500 are discarded.² Temporal H-MDT models are learned from fine to coarse scale. At the finer scale, there are 6×10 windows \mathcal{R}_i^1 on Ped1 (8×11 for Ped2), each covering a 41×41 pixel area and overlapping by 25 percent with each of its four neighbors. An MDT of five components is learned per window. At coarser spatial scales, an MDT is estimated from the MDTs of the four regions that it covers at the immediately finer resolution. Each estimated MDT has one more component than its ancestor MDTs. Overall, there are 10 scales in Ped1 and 11 in Ped2. Spatial anomaly maps use a 31×31 center window and surround windows of size equivalent to \mathcal{R}_i^s . For segmentation, $7 \times 7 \times 10$ patches are extracted from the 40 frames surrounding that under analysis. There are five DT components at all levels of the spatial hierarchy. Both temporal and spatial MDTs have an eight-dimensional state space. The sensitivity of the proposed detector to some of these parameters is discussed in Appendix C.2, available in the online supplemental material.

6.4 Descriptor Comparison

The first experiment evaluated the benefits of MDT-based over optical flow descriptors. The optical flow descriptors considered were the local motion histogram (LMH) of [19], the force flow descriptor of [4], and the mixture of optical flow models (MPPCA) of [3]. LMH uses statistics of local motion, and is representative of traditional background subtraction representations, force flow is a descriptor for spatial anomaly detection, and MPPCA a temporal anomaly detector. For the MDT, only the anomaly maps of finest temporal and coarsest spatial scale were considered here. Since the goal was to compare descriptors, the high-level components of the models in which they were proposed, for example, the LDA of [4], the MRF of [3], and the proposed CRF, were not used. Instead, anomaly predictions were smoothed with a simple

2. This variance threshold is quite conservative, only eliminating regions of very little motion. For the data sets used in our experiments, this has not led to the elimination of any objects from further consideration. In other contexts, for example, scenes where objects are static for periods of time, this could happen. In this case, the threshold should be set to zero.

TABLE 2
Descriptor Performance on UCSD Anomaly Data Set

Descriptor	Criterion			
	EER		RD	
	Ped1	Ped2	Ped1*	Ped2
MDT-temp.	22.9%	27.9%	59.3% (48.2%)	56.8%
MDT-spat.	43.8%	28.7%	50.8% (54.2%)	63.4%
MPPCA	35.6%	35.8%	23.2% (27.0%)	22.4%
force flow	36.5%	35.0%	40.9% (38.8%)	27.6%
LMH	38.9%	45.8%	32.6% (35.3%)	22.4%

*numbers outside/inside parentheses are results by full/partial annotation (same for the rest of the paper).

$20 \times 20 \times 10$ Gaussian filter. Anomaly predictions were generated by thresholding the filtered anomaly maps and ROC curves by varying thresholds.

The performance of the different descriptors, under both the frame-level (EER) and pixel-level (RD) criteria (using both full and partial annotation in Ped1), is summarized in Table 2. The corresponding ROC curves are presented in Appendix C.1 (Fig. 13), available in the online supplemental material. Examples of detected anomalies are shown in Fig. 6. Under the frame-level criterion, temporal MDT has the best performance in both scenes. Spatial MDT performs worse than others in Ped1 but ranks second in Ped2. However, for the more precise pixel-level criterion, spatial MDT is the top or second best performer. In this case, both MDTs significantly outperform all optical flow descriptors. The gap between corresponding competitors (e.g., temporal MDT versus MPPCA or LMH, spatial MDT versus force flow) is of at least 10 percent RD. These results show that there is a definite benefit to the joint representation of appearance and dynamics of the MDT.

This is not totally surprising, given the limitations of optical flow. First, the brightness constancy assumption is easily violated in crowded scenes, where stochastic motion and occlusions prevail. Second, optical flow measures instantaneous displacement, while the DT is a smooth motion representation with extended temporal support. Finally, while optical flow is a bandpass measure, which eliminates most of the appearance information, the DT models both appearance and dynamics. The last two properties are particularly important for crowded scenes, where objects occlude and interact in complicated manners.

TABLE 3
Filter Performance on the UCSD Anomaly Data Set

Criterion	Scene	Method		
		S-MDT w/ smoothing	H-MDT w/ smoothing	H-MDT w/ CRF-filtering
Frame-Level (EER)	Ped1	21.3%	21.6%	17.8%
	Ped2	23.9%	22.3%	18.5%
Pixel-Level (RD)	Ped1	59.7% (56.7%)	71.2% (60.1%)	74.5% (64.8%)
	Ped2	65.0%	69.6%	70.1%

Overall, although optical flow can signal fast moving anomalous subjects, it leads to too many false positives in regions of complex motion, occlusion, and so on. More interesting is the lack of advantage for either spatial or temporal anomaly detection, both among MDT maps and prior techniques (no clear advantage to either force flow or MPPCA). In fact, as shown in Fig. 6, temporal and spatial anomalies tend to be different objects. This suggests the combination of the two strategies.

6.5 Scale and Globally Consistent Prediction

We next investigated the benefits of information fusion across space and scale, with the proposed CRF. We started with a *single-scale* description (S-MDT), using only the anomaly maps at finest temporal and coarsest spatial scales, i.e., a 3D feature per site. We next considered a *multiscale* description, using the whole H-MDT. In both cases, inference was performed with logistic regression, i.e., the interaction term of (16) turned off, and the Gaussian filter of the previous section. In each trial, the logistic classifier was trained by Newton’s method [42]. Finally, we considered the full blown CRF, denoted *CRF filter*. The dimensions of the spatial and temporal CRF neighborhoods were set to $|\mathcal{N}^S| = 6$, $|\mathcal{N}^T| = 3$. ROC curves were generated by varying the threshold for prediction.

Table 3 presents a comparison of the three approaches. The corresponding ROC curves are shown in Appendix C.1 (Fig. 14), available in the online supplemental material. Under the pixel-level criterion, the multiscale maps have higher accuracy than their single-scale counterparts, demonstrating the benefits of modeling anomalies in scale space (improvement of RD by as much as 11 percent). The CRF



Fig. 6. Anomaly predictions of temporal MDT, spatial MDT, MPPCA, force flow, and LMH (from left to right). Red regions are abnormal pixels. All predictions generated with thresholds such that the different approaches have similar FPR under frame-level protocol (these settings apply to all the subsequent figures unless otherwise stated).



Fig. 7. Examples of anomaly localization with Gaussian smoothing (in blue) and CRF filter (in red). The latter predicts more accurately the spatiotemporal support of anomalies in crowded regions, where occlusion is prevalent.

TABLE 4
Performance of Various Methods (RD/Seconds per Frame) by Pixel-Level Criterion on UCSD Anomaly Data Set

	MDT (temporal)	MDT (spatial)	H-MDT-temp. w/ CRF-filtering	H-MDT-spat. w/ CRF-filtering	H-MDT w/ CRF-filtering [◇]	sparse recon- struction [8] [‡]	Bayesian video parsing [9] [#]
Ped1	- / 0.52	- / 0.57	65% (52%) / 0.61	57% (58%) / 0.65	75% (65%) / 1.11	-(46%) / 3.8	77% (68%) / 5~10
Ped2	- / 0.64	- / 0.69	55% / 0.76	68% / 0.80	70% / 1.38	-	-

Implementation: [◇] C/2.8-GHz CPU/2-GB RAM; [‡] C++ and Matlab (feature extraction and model inference)/2.6GHz CPU/2GB RAM; [#]Matlab/dual-core 2.7GHz CPU/8GB RAM.

filter further improves performance (improvement of RD by as much as 3 percent), demonstrating the gains of globally consistent inference. As shown in Fig. 7, the visual improvements are even more substantial.³ Simple filtering does not take into account interactions between neighboring sites and smooths the anomaly maps uniformly. On the other hand, the CRF adapts the degree of smoothing to the spatiotemporal structure of the anomalies, increasing the precision of anomaly localization. Note how, in Fig. 7, the CRF-filter successfully excludes occluded but normally behaving pedestrians from anomaly regions. These improvements are not always captured by the frame-level criterion. In fact, there is little EER difference between S-MDT and H-MDT. The inconsistency between frame- and pixel-level results in Tables 2 and 3 shows that the former is not a good measure of anomaly detection performance. Henceforth, only the pixel-level criterion is used in the remaining experiments on this data set.

6.6 Anomaly Detection Performance

We next evaluated the performance of the complete anomaly detector. For this, we selected two detectors from the recent literature, with state-of-the-art performance for temporal [8] and combined spatial and temporal anomaly detection [9]. The RD of the various methods is summarized in Table 4, for both partial and full annotation. The corresponding ROC curves are shown in Fig. 8. Table 4 also presents the processing time per video frame of each method. Missing entries indicate unavailable results for the particular data set and/or annotation type. A discussion of the detection errors made by the detector is given in Appendix C.3, available in the online supplemental material.

On Ped1, the temporal component of the proposed detector substantially outperforms the temporal detector of [8]. A multiple-scale temporal anomaly map with CRF filtering increases the 46 percent RD⁴ of [8] to 52 percent. A similar implementation of the spatial anomaly detector (a multiple-scale map plus CRF filtering) achieves 58 percent. Combining both maps and multiple spatial scales further

improves the RD to 65 percent. Computationally, the proposed detector is also much more efficient. For implementations on similar hardware (see footnotes of Table 4), it requires 1.11 s/frame, as compared to the 3.8 s/frame reported for [8].

Like the proposed detector, the Bayesian video parsing (BVP) of [9] combines spatial and temporal anomaly detection, using a more complex video representation, parsing of the video to extract all the objects in the scene, a support vector machine classifier for detection of temporal anomalies, a graphical model with seven nodes per site (and multiple nonparametric models for location, scale, and velocity) for detection of spatial anomalies, and occlusion reasoning. This is an elegant solution, which achieves slightly better RD than the proposed detector (2 percent for full and 3 percent for partial annotation), but at substantially higher computational cost (5 to 10 times slower). We believe that when both accuracy and computation are considered, the proposed detector is a more effective solution. However, these results suggest that gains could be achieved by expanding the proposed CRF, as [9] trades a much simpler representation of video dynamics (optical flow *versus* MDT) for more sophisticated inference. It would be interesting to consider CRF extensions with some of the properties of the graphical model of [9], namely, explicit occlusion reasoning. This is left for subsequent research.

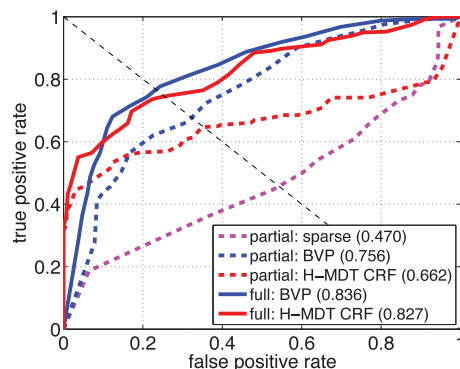


Fig. 8. ROC curves of pixel-level criterion on Ped1.

3. More results at <http://www.svl.ucsd.edu/projects/anomaly/results.html>.

4. These numbers refer to partial annotation, the only available for [8].

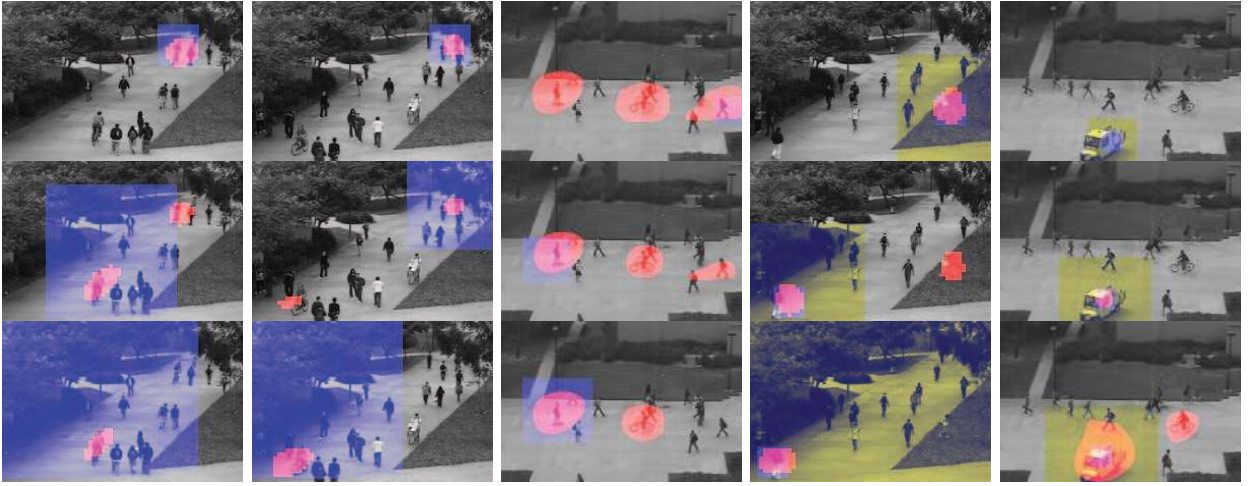


Fig. 9. Impact of context on anomaly maps. First three columns: Temporal anomalies, cell coverage at different HMDT layers shown in blue. Last two columns: Spatial anomalies, example center (surround) windows shown in blue (light yellow).

6.7 Role of Context in Anomaly Judgments

We next investigated the impact of normalcy context in anomaly judgments. For temporal anomalies, context is determined by the subregion size: As the latter increases, temporal models become more global. Fig. 9 shows that the scale of normalcy context significantly impacts anomaly scores. For example, the two cyclists on the left-most columns of the figure are missed at small scales but detected by the more global models. On the other hand, a leftward heading pedestrian in the third column has high anomaly score at the finest scale but is not anomalous in larger contexts. In summary, no single context is effective for all scenes. Due to the stochastic arrangements of people within crowds, two crowds of the same size can require different context sizes. In general, the optimal size depends on the crowd configuration and the anomalous event.

A similar observation holds for spatial anomalies, where context is set by the size of the surround window. For example, in the fourth column of Fig. 9, the subject walking on the grass is very salient when compared to her immediate neighbors, and anomaly detection benefits from a narrower context. For larger contexts, she becomes less unique than a man that walks in the direction opposite to his neighbors. On the other hand, the cart and bike of the last column only pop out when the surround window is large enough to cover some pedestrians. In summary, anomalies depend strongly on scene context, and this dependence can vary substantially from scene to scene. It is, thus, important to fuse anomaly information across spatial scales.

6.8 Performance on Other Benchmark Data Sets

The detection of anomalous events in crowded scenes can be evaluated in a few data sets other than UCSD. These have various limitations in terms of size, saliency of the anomalies, evaluation criteria, and so on. They are discussed in this section where, for completeness, we also present the results of the proposed anomaly detector.

UMN. The UMN data set⁵ contains three escape scenes. Normal events depict individuals wandering around or organized in groups. Abnormal events depict a crowd escaping in panic. Each scene contains several normal-

abnormal events (e.g., seconds of normalcy followed by a short abnormal event). The main limitations of this data set are that

1. it is relatively small (scenes 1, 2, and 3 contain two, six, and three anomaly instances),
2. it has no pixel-level ground truth,
3. the anomalies are staged, and
4. it produces very salient changes in the average motion intensity of the scene.

As a result, several methods achieve near perfect detection.

The proposed detector was based on 3×3 subregions of size 180×180 at the finest spatial scale and a 3-scale anomaly map for both the temporal and spatial components. One normal-abnormal instance of each scene was used to train the temporal normalcy model and CRF filter, and the remaining instances for testing. A comparison to previous results in the literature, under the frame-level criterion, is presented in Table 5 and Fig. 11. Due to the salient motion discontinuities, the temporal component (99.2 percent AUC) substantially outperforms the spatial component (97.9 percent). Nevertheless, the complete detector achieves the best performance (99.5 percent). This is nearly perfect, and comparable to the previous best results in the literature.

Subway. The Subway data set [19] consists of two sequences recorded from the entrance (1 h and 36 min, 144,249 frames) and exit (43 min, 64,900 frames) of a

TABLE 5
Anomaly Detection Performance in AUC/ERR (Percent)

method	UMN	Subway		U-turn
		entrance	exit	
chaotic invariants [1]	99.4/5.3	-/-	-/-	-/-
social force [4]	94.9/12.6	-/-	-/-	-/-
sparse [8]	99.6/2.8	80.2/26.4	83.3/24.4	-/-
local stat. aggr. [23]	99.5/3.4	88.4/17.9	-/-	94.7/-
H-MDT-spat. CRF	97.9/7.8	68.2/37.0	67.0/34.1	83.9/-
H-MDT-temp. CRF	99.2/5.3	88.9/18.6	87.5/17.9	92.9/-
H-MDT CRF	99.5/3.7	89.7/16.4	90.8/16.7	95.2/-

5. <http://mha.cs.umn.edu/Movies/Crowd-Activity-All.avi>.



Fig. 10. Anomalies detected by H-MDT CRF on the UMN (left), Subway (center), and U-turn (right) data sets.

subway station. Normal behaviors include people entering and exiting the station; abnormal consist of people moving in the wrong direction (exiting the entrance or entering the exit) or avoiding payment. The main limitations of this data set are: 1) reduced number of anomalies, and 2) predictable spatial localization (entrance and exit regions). The original 512×384 frames were down sampled to 320×240 , and 2×3 subregions of size 90×90 , covering either the entrance or exit regions, were used at the finest spatial scale. A 3-scale anomaly map was computed for both spatial and temporal anomalies. Video patches were of size $15 \times 15 \times 15$, and 10 min of video from each sequence was used to train the temporal normalcy model and CRF filters, while the remaining video was used for testing. Table 5 and Fig. 11 present a comparison of the proposed detector against recently published results on this data set. Again, the temporal component outperforms its spatial counterpart, but the best performance is obtained by combination of both temporal and spatial anomaly maps (H-MDT CRF). This achieves the best result among all methods, outperforming the sparse reconstruction of [8] and the local statistical aggregates of [23]. Note that, for this data set, the gains in both AUC and EER are substantial.

U-turn. The U-turn data set [5] consists of one video sequence (roughly 6,000 frames of size 360×240) recorded by a static camera overlooking the traffic at a road intersection. The video is split into two clips of equal length for cross validation and anomalies consist of illegal vehicle motion at the intersection. The main limitations of this data set are: 1) the limited size, 2) absence of pixel-level ground truth, and 3) sparseness of the scenes. The latter enables the use of object-based operations, for example, tracking and analysis of object trajectories [5], which we do not exploit.

For temporal anomaly detection, MDTs were learned using $20 \times 20 \times 30$ patches from 3×4 subregions covering the intersection. This was the finest level of a 3-scale hierarchical model. For spatial anomaly detection, segmentation was computed with a 5-component MDT learned from $15 \times 15 \times 30$ patches extracted from 45 consecutive frames. An observation lattice of step $15 \times 15 \times 10$ was used to evaluate anomaly scores, and the neighborhood size of the CRF filter was 2. The performance of the detector is summarized in Table 5 and Fig. 11. Due to the sparsity of the scenes (not enough spatial context around cars making illegal turns to establish them as anomalous) the performance of the spatial anomaly detector is quite weak. However, the combination of the spatial and temporal anomaly maps again outperforms the temporal channel, achieving the best performance. Overall, the proposed detector has the best AUC on this data set. Examples of detected anomalies, for this and the other two data sets, are shown in Fig. 10.

7 CONCLUSION

In this work, we proposed an anomaly detector that spans time, space, and spatial scale, using a joint representation of video appearance and dynamics and globally consistent inference. For this, we modeled crowded scenes with a hierarchy of MDT models, equated temporal anomalies to background subtraction, spatial anomalies to discriminant saliency, and integrated anomaly scores across time, space, and scale with a CRF. It was shown that the MDT representation substantially outperforms classical optical flow descriptors, that spatial and temporal anomaly detection are complementary processes, that there is a benefit to defining anomalies with respect to various normalcy contexts, i.e., in anomaly scale space, and that it is important to guarantee globally consistent inference across space, time and scale. We have also introduced a challenging anomaly detection data set, composed of complex scenes of pedestrian crowds, involving stochastic motion, complex occlusions, and object interactions. This data set provides both frame-level and pixel-level ground truth, and a protocol for the evaluation of anomaly detection algorithms. The proposed anomaly detector was shown effective on both this and a number of previous data sets. When compared to previous methods, it outperformed various state-of-the-art approaches, either in absolute performance or in terms of the tradeoff between anomaly detection accuracy and complexity.

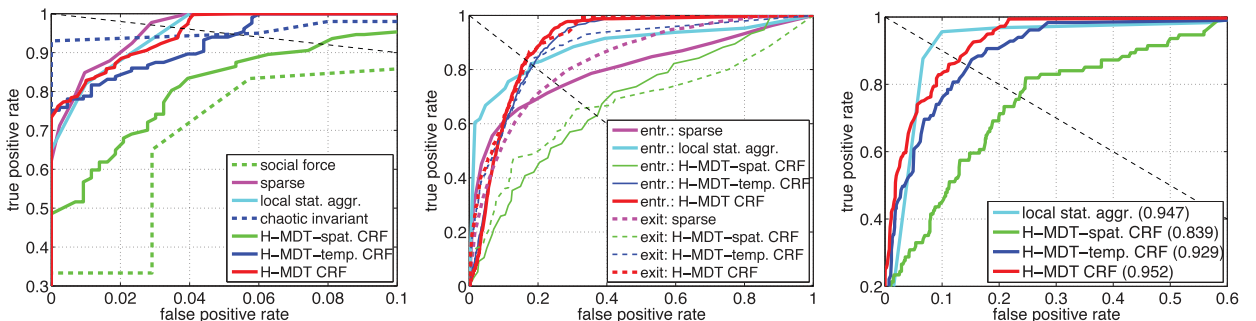


Fig. 11. ROC curves of frame-level criterion on the UMN (left), Subway (center), and U-turn (right) data sets.

REFERENCES

- [1] S. Wu, B. Moore, and M. Shah, "Chaotic Invariants of Lagrangian Particle Trajectories for Anomaly Detection in Crowded Scenes," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [2] L. Kratz and K. Nishino, "Anomaly Detection in Extremely Crowded Scenes Using Spatio-Temporal Motion Pattern Models," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [3] J. Kim and K. Grauman, "Observe Locally, Infer Globally: A Space-Time MRF for Detecting Abnormal Activities with Incremental Updates," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [4] R. Mehran, A. Oyama, and M. Shah, "Abnormal Crowd Behavior Detection Using Social Force Model," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [5] Y. Benezeth, P. Jodoin, V. Saligrama, and C. Rosenberger, "Abnormal Events Detection Based on Spatio-Temporal Co-Occurrences," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [6] A. Basharat, A. Gritai, and M. Shah, "Learning Object Motion Patterns for Anomaly Detection and Improved Object Detection," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2008.
- [7] T. Xiang and S. Gong, "Video Behavior Profiling for Anomaly Detection," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, no. 5, pp. 893-908, May 2008.
- [8] Y. Cong, J. Yuan, and J. Liu, "Sparse Reconstruction Cost for Abnormal Event Detection," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [9] B. Antić and B. Ommer, "Video Parsing for Abnormality Detection," *Proc. IEEE Int'l Conf. Computer Vision*, 2011.
- [10] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly Detection: A Survey," *ACM Computing Surveys*, vol. 41, no. 3, article 15, 2009.
- [11] G. Doretto, A. Chiuso, Y. Wu, and S. Soatto, "Dynamic Textures," *Int'l J. Computer Vision*, vol. 51, no. 2, pp. 91-109, 2003.
- [12] A. Chan and N. Vasconcelos, "Modeling, Clustering, and Segmenting Video with Mixtures of Dynamic Textures," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, no. 5, pp. 909-926, May 2008.
- [13] D. Gao and N. Vasconcelos, "Decision-Theoretic Saliency: Computational Principles, Biological Plausibility, and Implications for Neurophysiology and Psychophysics," *Neural Computation*, vol. 21, no. 1, pp. 239-271, 2009.
- [14] C. Stauffer and W. Grimson, "Learning Patterns of Activity Using Real-Time Tracking," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 747-757, Aug. 2000.
- [15] T. Zhang, H. Lu, and S. Li, "Learning Semantic Scene Models by Object Classification and Trajectory Clustering," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009.
- [16] N. Siebel and S. Maybank, "Fusion of Multiple Tracking Algorithms for Robust People Tracking," *Proc. European Conf. Computer Vision*, 2006.
- [17] X. Cui, Q. Liu, M. Gao, and D.N. Metaxas, "Abnormal Detection Using Interaction Energy Potentials," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [18] F. Jiang, J. Yuan, S.A. Tsafaris, and A.K. Katsaggelos, "Anomalous Video Event Detection Using Spatiotemporal Context," *Computer Vision and Image Understanding*, vol. 115, no. 3, pp. 323-333, 2011.
- [19] A. Adam, E. Rivlin, I. Shimshoni, and D. Reinitz, "Robust Real-Time Unusual Event Detection Using Multiple Fixed-Location Monitors," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 30, no. 3, pp. 555-560, Mar. 2008.
- [20] B. Zhao, L. Fei-Fei, and E. Xing, "Online Detection of Unusual Events in Videos via Dynamic Sparse Coding," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2011.
- [21] D. Helbing and P. Molnár, "Social Force Model for Pedestrian Dynamics," *Physical Rev. E*, vol. 51, no. 5, pp. 4282-4286, 1995.
- [22] O. Boiman and M. Irani, "Detecting Irregularities in Images and in Video," *Int'l J. Computer Vision*, vol. 74, no. 1, pp. 17-31, 2007.
- [23] V. Saligrama and Z. Chen, "Video Anomaly Detection Based on Local Statistical Aggregates," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2012.
- [24] R. Hamid, A. Johnson, S. Batta, A. Bobick, C. Isbell, and G. Coleman, "Detection and Explanation of Anomalous Activities: Representing Activities as Bags of Event N-Grams," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2005.
- [25] D. Zhang, D. Gatica-Perez, S. Bengio, and I. McCowan, "Semi-Supervised Adapted HMMs for Unusual Event Detection," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2005.
- [26] C. Stauffer and W. Grimson, "Adaptive Background Mixture Models for Real-Time Tracking," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 1999.
- [27] L. Itti, C. Koch, and E. Niebur, "A Model of Saliency-Based Visual Attention for Rapid Scene Analysis," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 11, pp. 1254-1259, Nov. 1998.
- [28] R. Shumway and D. Stoffer, "An Approach to Time Series Smoothing and Forecasting Using the EM Algorithm," *J. Time Series Analysis*, vol. 3, no. 4, pp. 253-264, 1982.
- [29] S. Roweis and Z. Ghahramani, "A Unifying Review of Linear Gaussian Models," *Neural Computation*, vol. 11, no. 2, pp. 305-345, 1999.
- [30] S. Kullback, *Information Theory and Statistics*. Dover Publications, 1968.
- [31] D. Gao, V. Mahadevan, and N. Vasconcelos, "On the Plausibility of the Discriminant Center-Surround Hypothesis for Visual Saliency," *J. Vision*, vol. 8, no. 7, pp. 1-18, 2008.
- [32] V. Mahadevan and N. Vasconcelos, "Background Subtraction in Highly Dynamic Scenes," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2008.
- [33] A. Chan and N. Vasconcelos, "Probabilistic Kernels for the Classification of Auto-Regressive Visual Processes," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2005.
- [34] J.R. Hershey and P.A. Olsen, "Approximating the Kullback Leibler Divergence between Gaussian Mixture Models," *Proc. IEEE Int'l Conf. Acoustics, Speech, and Signal Processing*, 2007.
- [35] A.B. Chan and N. Vasconcelos, "Efficient Computation of the KL Divergence between Dynamic Textures," Technical Report SVCL-TR-2004-02, Dept. of Electrical and Computer Eng., Univ. of California San Diego, 2004.
- [36] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2006.
- [37] A. Chan, E. Coviello, and G. Lanckriet, "Clustering Dynamic Textures with the Hierarchical EM Algorithm," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [38] J. Lafferty, A. McCallum, and F. Pereira, "Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data," *Proc. 18th Int'l Conf. Machine Learning*, 2001.
- [39] S. Kumar and M. Hebert, "Discriminative Fields for Modeling Spatial Dependencies in Natural Images," *Proc. Advances in Neural Information Processing Systems*, 2004.
- [40] X. He, R. Zemel, and M. Carreira-Perpinán, "Multiscale Conditional Random Fields for Image Labeling," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2004.
- [41] G.E. Hinton, "Training Products of Experts by Minimizing Contrastive Divergence," *Neural Computation*, vol. 14, pp. 1771-1800, 2002.
- [42] T. Minka, "A Comparison of Numerical Optimizers for Logistic Regression," technical report, Microsoft Research, 2003.
- [43] V. Mahadevan, W. Li, V. Bhalodia, and N. Vasconcelos, "Anomaly Detection in Crowded Scenes," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2010.
- [44] T.P. Kah-Kay Yung, "Example-Based Learning for View-Based Human Face Detection," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 20, no. 1, pp. 39-51, Jan. 1998.



Weixin Li received the bachelor's degree from Tsinghua University, Beijing, China, in 2008, the MSc degree in electrical engineering from the University of California, San Diego, in 2011, and is currently working toward the PhD degree. His research interests primarily include computational vision and machine learning, with specific focus on visual analysis of human behavior, activity, and event, and models with latent variables and their applications. He is a student member of the IEEE.



Vijay Mahadevan received the BTech degree from the Indian Institute of Technology, Madras, in 2002, the MS degree from Rensselaer Polytechnic Institute, Troy, New York, in 2003, and the PhD degree from the University of California, San Diego, in 2011, all in electrical engineering. From 2004 to 2006, he was with the Multimedia group at Qualcomm Inc., San Diego, California. He is currently with Yahoo! Labs, Bengaluru. His interests include computer

vision and machine learning and their applications. He is a member of the IEEE.



Nuno Vasconcelos received the licenciatura degree in electrical engineering and computer science from the Universidade do Porto, Portugal, and the MS and PhD degrees from the Massachusetts Institute of Technology. He is a professor in the Electrical and Computer Engineering Department, University of California, San Diego, where he heads the Statistical Visual Computing Laboratory. He has received a US National Science Foundation (NSF) CAREER

award, a Hellman Fellowship, and has authored more than 150 peer-reviewed publications. He is a senior member of the IEEE.

▷ **For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.**